

Fundamentos de Biologia Molecular

Curso de Licenciatura em Biologia
2º Ano, 1º Semestre
Ano Letivo 2018/2019

Componente Teórico-Prática



Ciências
ULisboa

Faculdade
de Ciências
da Universidade
de Lisboa

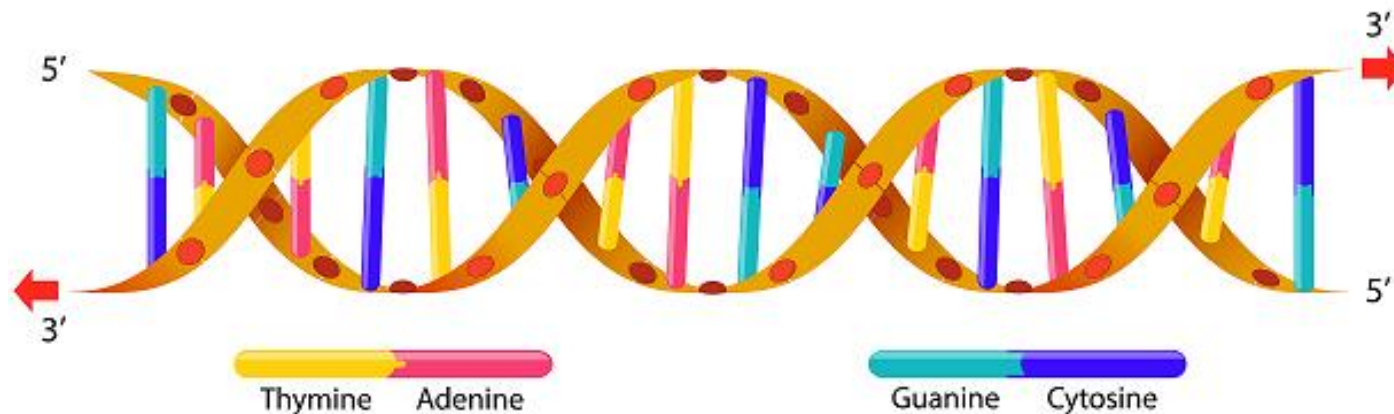
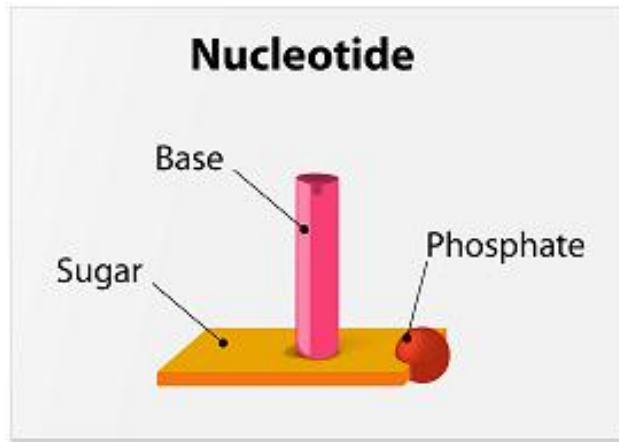
Docente Responsável: Rita Zilhão
Docente TPs: Andreia Figueiredo

TP3: DNA Sequencing

- Basic principle of Sanger sequencing: DNA structure
- Nobel prizes
- Dideoxy-terminating DNA/Sanger sequencing
- Overview of Sanger sequencing steps
- Technical advances
- Sequencing data analysis

TP3: DNA Sequencing

DNA structure



TP2: Polymerase Chain Reaction (PCR)



Frederick Sanger

▪ After his Ph.D. in 1943, Sanger started working for A. C. Chibnall, on identifying the free amino groups in insulin. In the course of identifying the amino groups, Sanger figured out ways to order the amino acids. He was the first person to obtain a protein sequence. By doing so, Sanger proved that proteins were ordered molecules and by analogy, the genes and DNA that make these proteins should have an order or sequence as well – first nobel prize in 1958

▪ Solving the problem of DNA sequencing became a natural extension of his work in protein sequencing. Sanger initially investigated ways to sequence RNA because it was smaller. Eventually, this led to techniques that were applicable to DNA and finally to the **dideoxy method most commonly used in sequencing reactions today**. Sanger won a second Nobel Prize for Chemistry in 1980 sharing it with Walter Gilbert, for their contributions concerning the determination of base sequences in nucleic acids, and Paul Berg for his work on recombinant DNA.

TP2: Polymerase Chain Reaction (PCR)



Paul Berg

■ An organism's genome is stored in the form of long rows of building blocks, known as nucleotides, which form DNA molecules. ***An organism's genome can be mapped by establishing the order of the nucleotides within the DNA molecule.*** In 1976, Allan Maxam and Walter ***Gilbert developed a method by which the ends of the DNA molecule could be marked using radioactive substances.*** After undergoing treatment with small amounts of chemicals that react with specific nucleotides, DNA fragments of varying lengths can be obtained. After undergoing what is known as electrophoresis, the nucleotide sequences in a DNA sample can be identified.



Walter Gilbert

■ DNA carries organisms' genomes and also determines their vital processes. The ability to artificially manipulate DNA opens the way to creating organisms with new characteristics. In conjunction with his studies of the tumor virus SV40, in 1972, Paul Berg succeeded in inserting DNA from a bacterium into the virus' DNA. He thereby **created the first DNA molecule made of parts from different organisms-"hybrid DNA" or "recombinant DNA"**.

TP3: DNA Sequencing



Frederick Sanger



Paul Berg

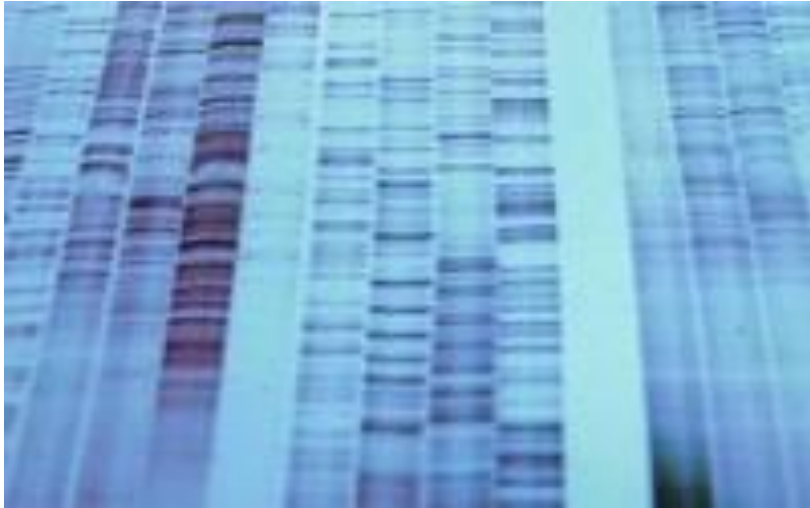


Walter Gilbert

The Nobel Prize in Chemistry 1980 was divided, one half awarded to Paul Berg "***for his fundamental studies of the biochemistry of nucleic acids, with particular regard to recombinant-DNA***", the other half jointly to Walter Gilbert and Frederick Sanger "***for their contributions concerning the determination of base sequences in nucleic acids***".

TP3: DNA Sequencing

Dideoxy-terminating DNA/Sanger sequencing concept



This method begins with the use of **special enzymes to synthesize fragments of DNA that terminate when a selected base appears in the stretch of DNA being sequenced.** These fragments are then sorted according to size by electrophoresis. Because of DNA's negative charge, the fragments move across the gel toward the positive electrode. The shorter the fragment, the faster it moves. Typically, each of the terminating bases within the collection of fragments is tagged with a **radioactive probe** for identification.

TP3: DNA Sequencing

Overview of Sanger sequencing steps

- first denature DNA – separation of double chain
- Anneal the primer (1 primer that anneals to the region of interest)
- The DNA is placed into 4 different tubes, one for each nitrogenous base
- DNA polymerase and 4 deoxynucleotides are added to each tube (dNTPs)
- One type of dideoxynucleotides is added to each tube
- DNA polymerase extends the DNA sequence (from de primer 5'-3')
- No nucleotide can be added to the DNA chain once a dideoxynucleotides has been incorporated, so each fragment will end with a labeled nucleotide.
- The content of each tube is denatured and separated by size by electrophoresis (polyacrylamide gel)
- So many sequences are synthesized that ddNTPs incorporation occurs in every sequence of the newly synthesized DNA sequence
- The further a specific strand has moved, the shorter it is – thus the position of the nucleotide that terminates that sequence can be determined by the distance travelled
- The order of nucleotides produced is a sequence (5'-3') that complements the original strand of DNA

TP3: DNA Sequencing

Overview of Sanger sequencing steps

Technique

DNA
(template strand)

5' C
T
G
A
C
T
T
C
G
A
C
A
A
3'

Primer

3'
T
G
T
T
5'

DNA
polymerase



TP3: DNA Sequencing

Overview of Sanger sequencing steps

Technique

DNA
(template strand)

5' C
T
G
A
C
T
T
C
G
A
C
A
A
3'

Primer

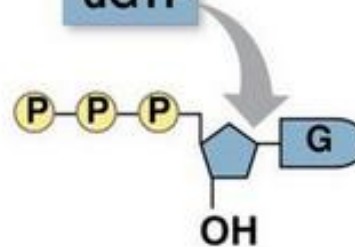
3'
T
G
T
T
5'

DNA
polymerase



Deoxyribo-
nucleotides

dATP
dCTP
dTTP
dGTP



TP3: DNA Sequencing

Overview of Sanger sequencing steps

Technique

DNA
(template strand)

5' CTGACTTCGACAA
3' A

Primer

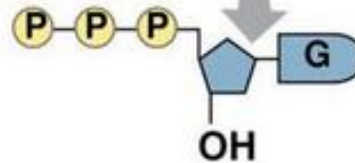
3' TGT
5' T

DNA
polymerase



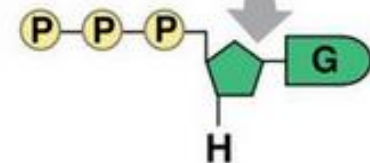
Deoxyribo-
nucleotides

dATP
dCTP
dTTP
dGTP



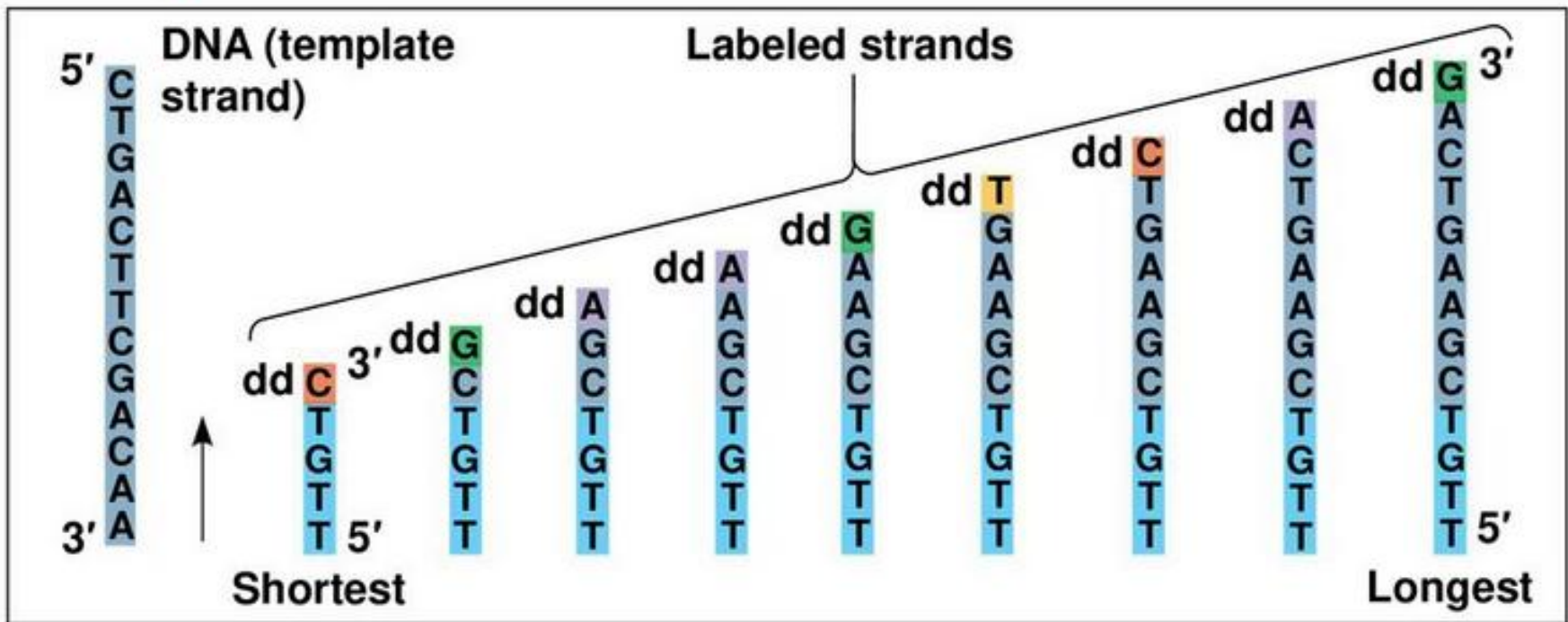
Dideoxynucleotides
(fluorescently tagged)

ddATP
ddCTP
ddTTP
ddGTP



TP3: DNA Sequencing

Technique



TP3: DNA Sequencing

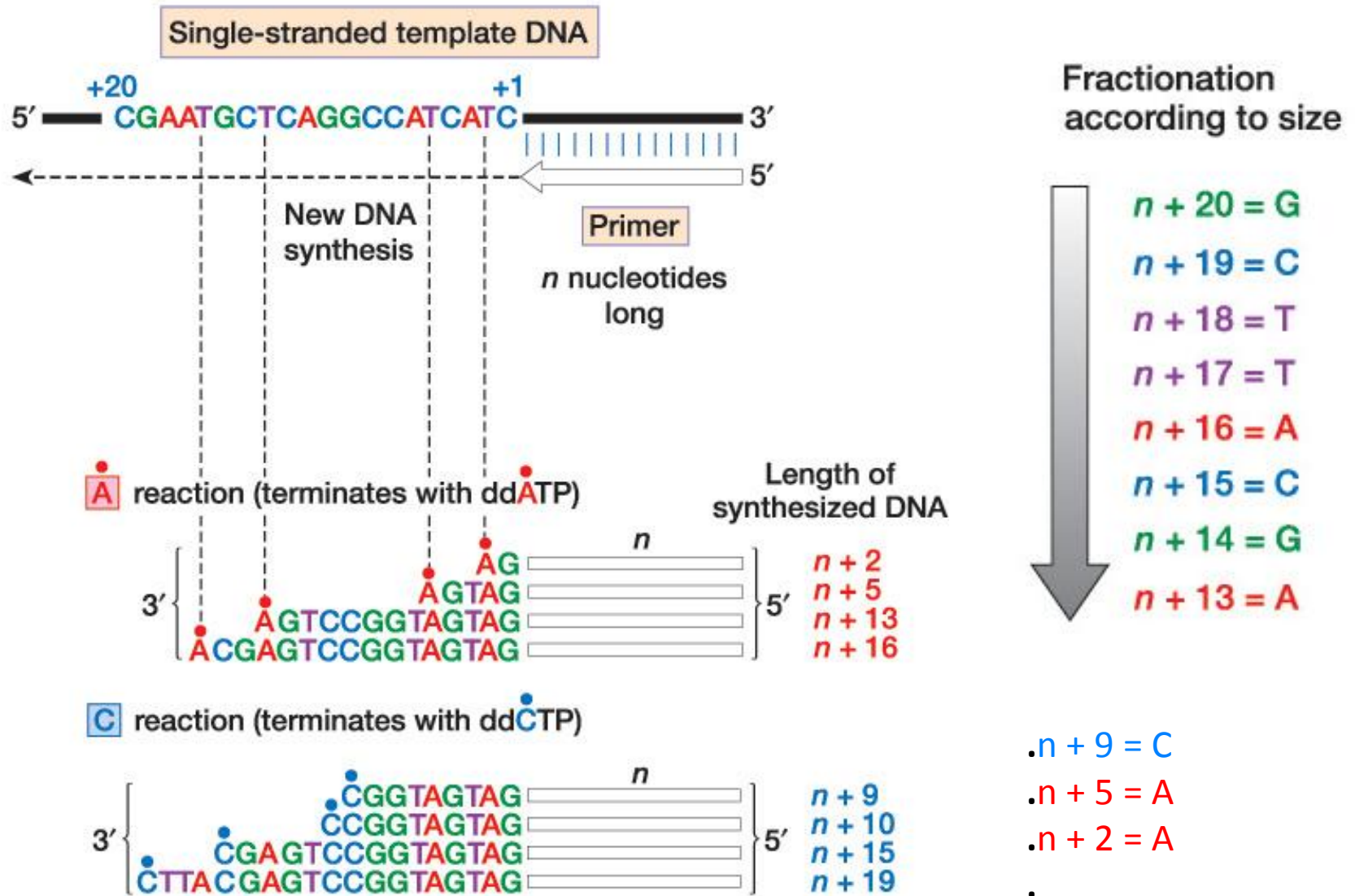
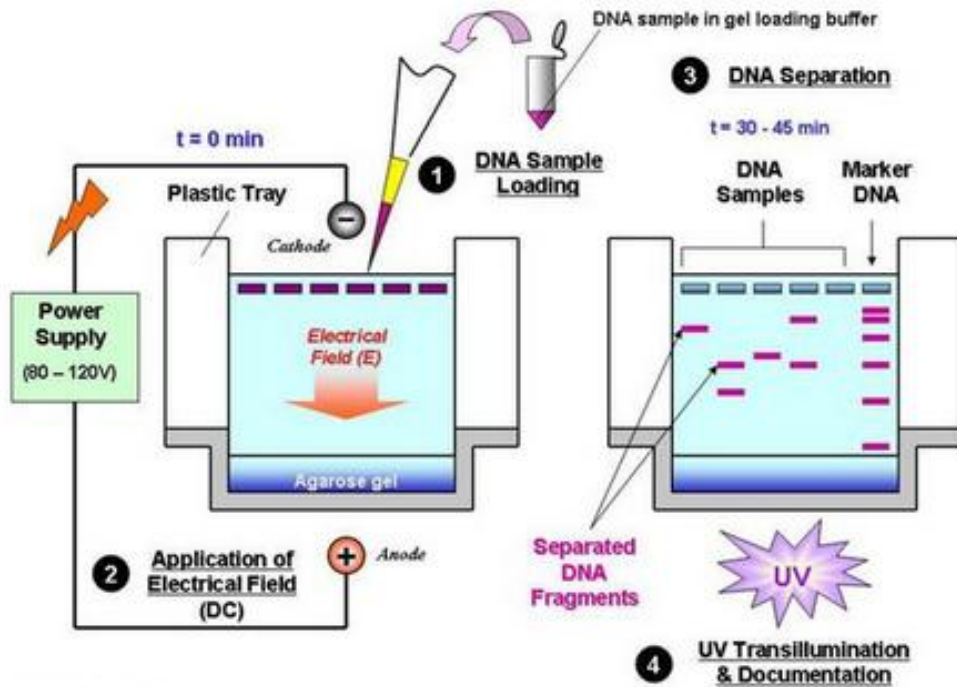


Figure 7-2 part 1 of 2 Human Molecular Genetics, 3/e. (© Garland Science 2004)

TP3: DNA Sequencing

Overview of Sanger sequencing steps



Graphic©ESchmid/2001

Polyacrylamide gel electrophoresis separates ssDNA molecules that differ in length by just **one nucleotide**

Molecules are labelled with a radioactive protein or radioactive isotope, visualized by autoradiography producing a banding pattern

<https://www.youtube.com/watch?v=3M0PyxFPwkQ>

<https://dnalc.cshl.edu/view/15479-Sanger-method-of-DNA-sequencing-3D-animation-with-narration.html>

TP3: DNA Sequencing

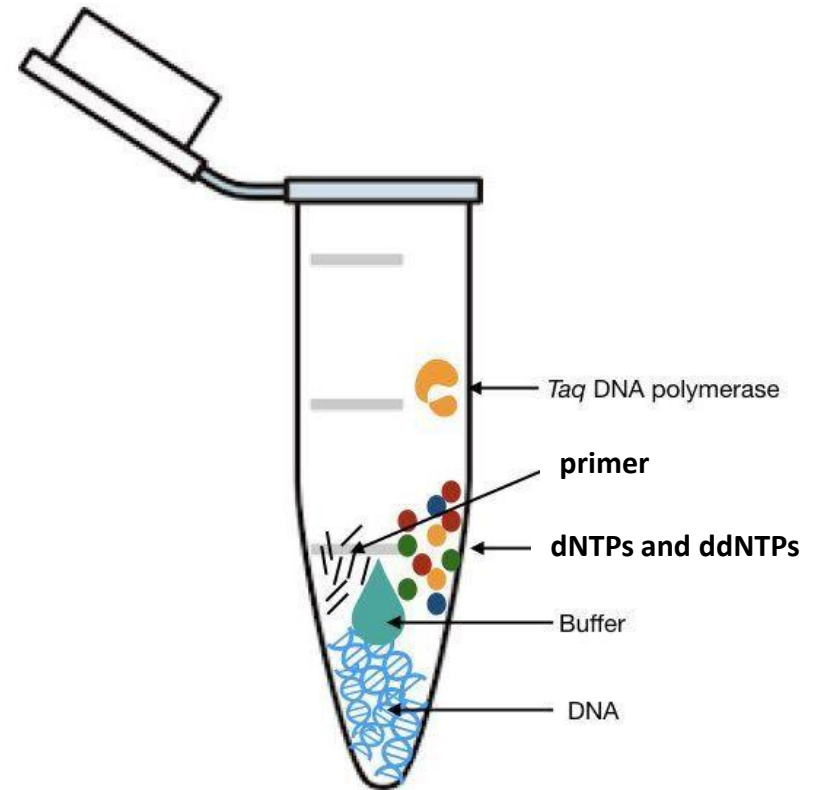
Sanger sequencing vs PCR

- **PCR is used to amplify DNA in its entirety.** While fragments of varying lengths may be produced by accident (e.g., the DNA polymerase might fall off), the goal is to duplicate the entire DNA sequence. To that end, the “ingredients” are the target DNA, nucleotides, DNA primer, and DNA polymerase (specifically Taq polymerase, which can survive the high temperatures required in PCR).
- **The goal of Sanger sequencing is to generate every possible length of DNA up to the full length of the target DNA.** That is why, in addition to the PCR starting materials, the **dideoxynucleotides** are necessary. Sanger sequencing and PCR can be brought together when generating the starting material for a Sanger sequencing protocol. PCR can be used to create many copies of the DNA that is to be sequenced. Having more than one template to work from makes the Sanger protocol more efficient.

TP3: DNA Sequencing

Dideoxy-terminating DNA sequencing reaction components

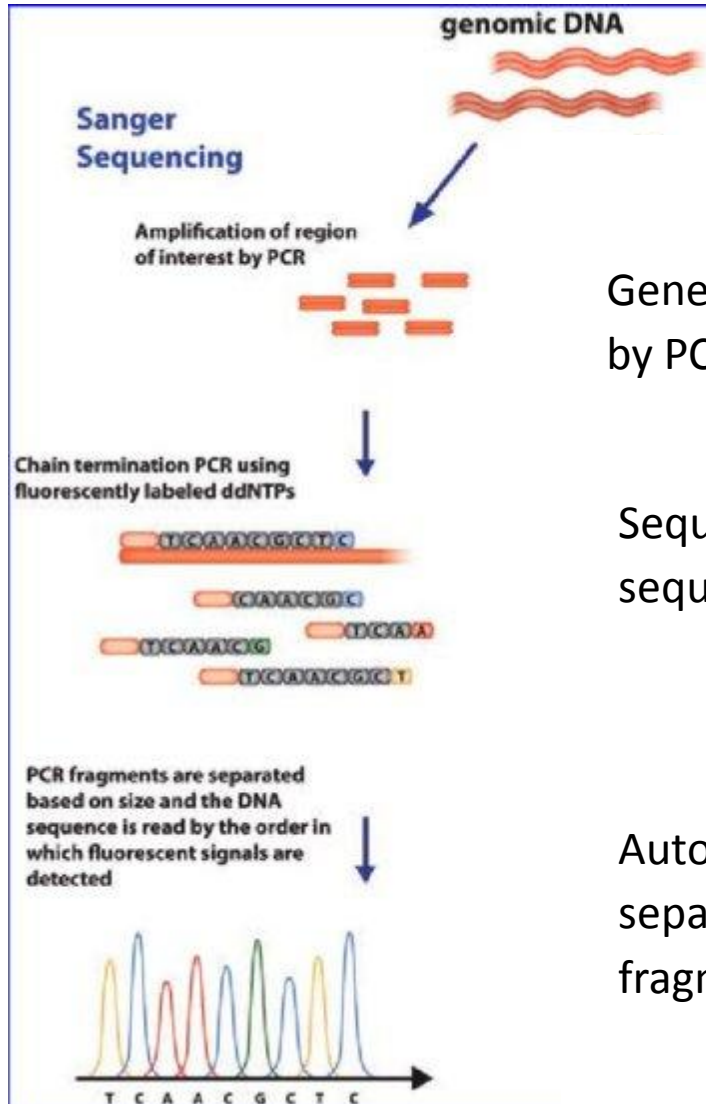
- **DNA template** to be sequenced
- One **specific primer** that binds to the template DNA and acts as a "starter" for the polymerase
- **nucleotides** (dATP, dTTP, dCTP, dGTP)
- **DNA polymerase** (proofreading activity, no 5'-3' exonuclease activity (eg Klenow fragment of E. coli polymerase, capacity of polymerizing ddNTPs , Eg. Vent)
- Dideoxy, or **chain-terminating**, versions of all four nucleotides (ddATP, ddTTP, ddCTP, ddGTP), each labeled (either radioactive label or fluorescent label with a different color of dye)



TP3: DNA Sequencing

DNA template to be sequenced

Specific DNA



Genomic DNA – billions of genes

Gene of interest – specifically amplified by PCR

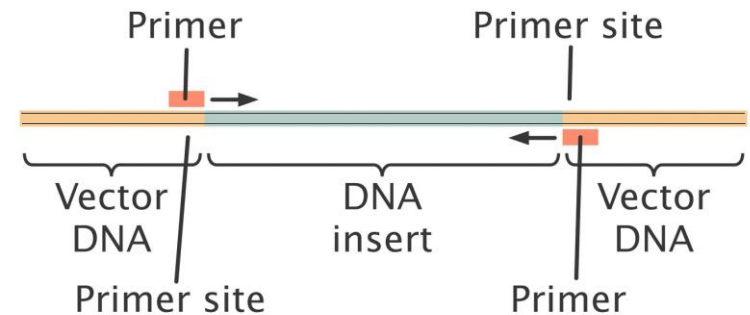
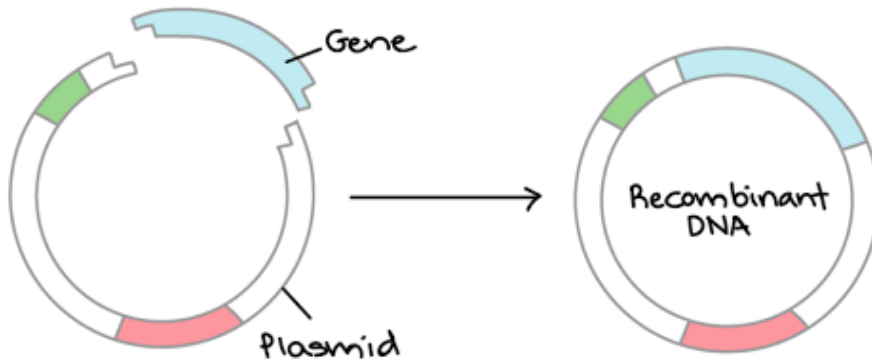
Sequencing reaction – Sanger sequencing (ddNTPs)

Automated sequencing by capillary separation of fluorescent labeled fragments

TP3: DNA Sequencing

DNA template to be sequenced

DNA cloned in a plasmid



A **universal sequencing primer** can be used to sequence many different template DNAs (eg M13, T7 primers)

Vectors contain it on either side of the site where DNA will be inserted

TP3: DNA Sequencing

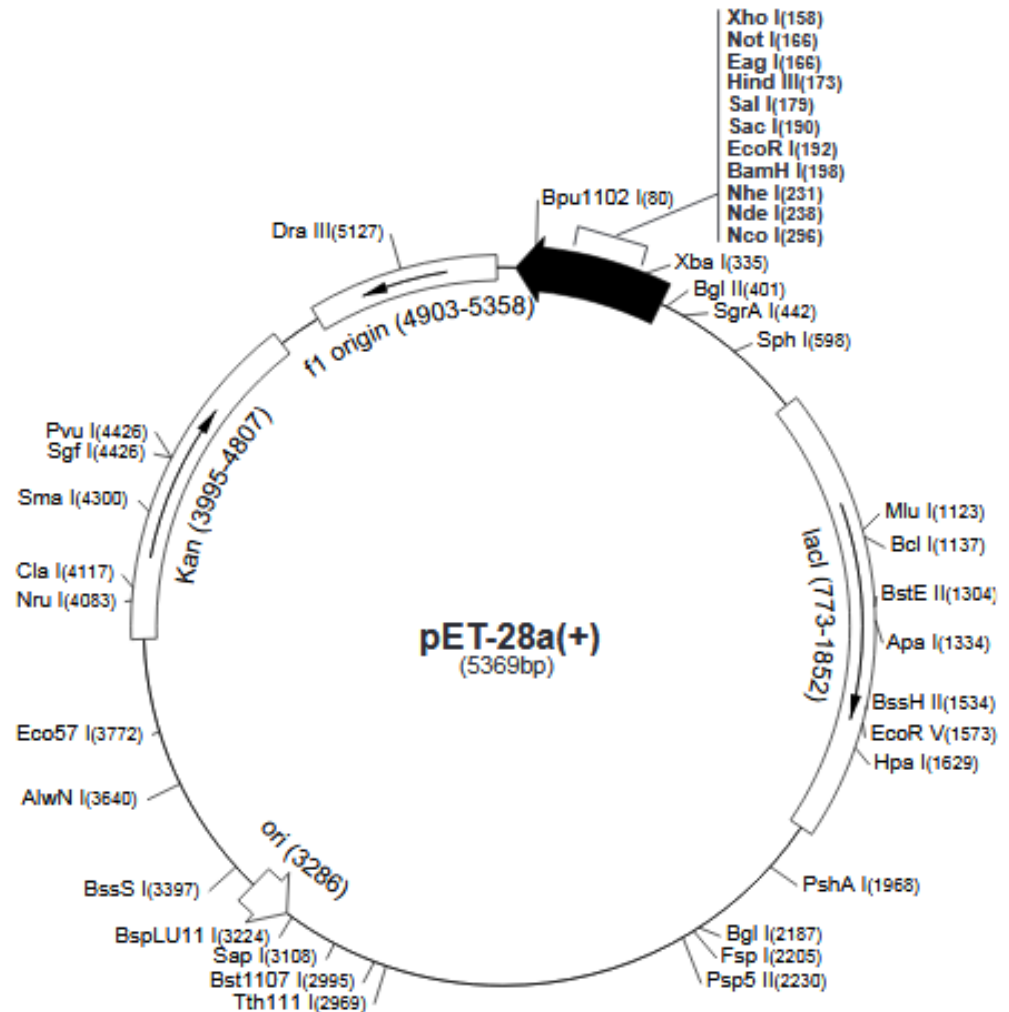
DNA template to be sequenced

DNA cloned in a plasmid

pET-28a(+) sequence landmarks

T7 promoter	370-386
T7 transcription start	369
His·Tag coding sequence	270-287
T7·Tag coding sequence	207-239
Multiple cloning sites (<i>Bam</i> H I - <i>Xho</i> I)	158-203
His·Tag coding sequence	140-157
T7 terminator	26-72
<i>lacI</i> coding sequence	773-1852
pBR322 origin	3286
Kan coding sequence	3995-4807
f1 origin	4903-5358

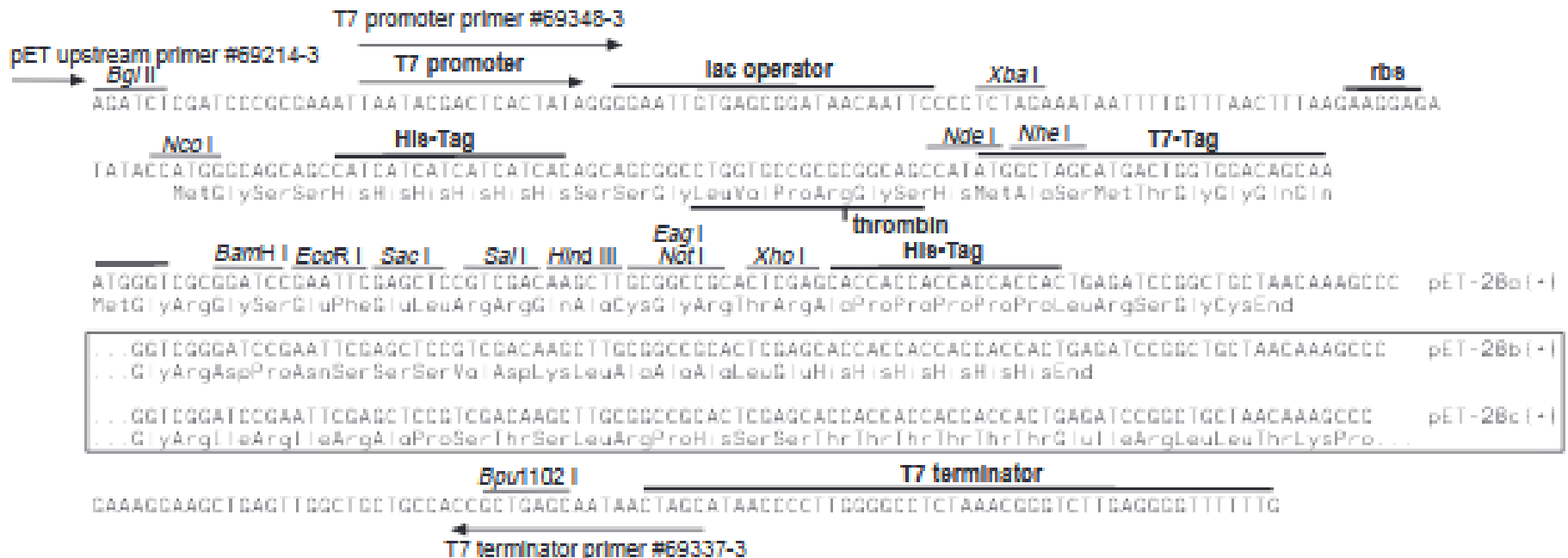
The maps for pET-28b(+) and pET-28c(+) are the same as pET-28a(+) (shown) with the following exceptions: pET-28b(+) is a 5368bp plasmid; subtract 1bp from each site beyond *Bam*H I at 198. pET-28c(+) is a 5367bp plasmid; subtract 2bp from each site beyond *Bam*H I at 198.



TP3: DNA Sequencing

DNA template to be sequenced

DNA cloned in a plasmid

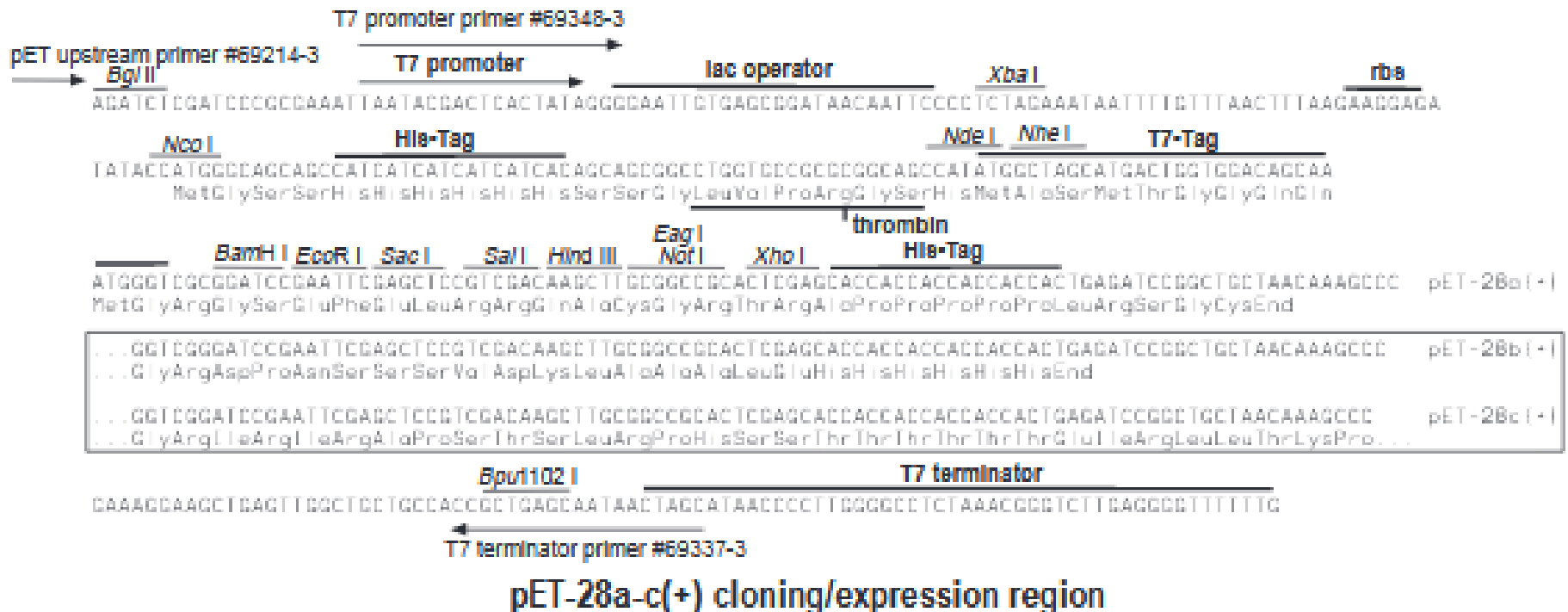


pET-28a-c(+) cloning/expression region

TP3: DNA Sequencing

DNA template to be sequenced

DNA cloned in a plasmid



TP3: DNA Sequencing

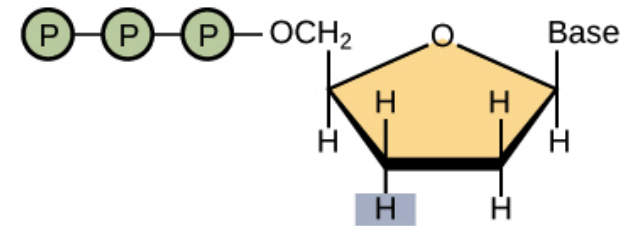
ddNTPs

- Dideoxy, or **chain-terminating**, versions of all four nucleotides (ddATP, ddTTP, ddCTP, ddGTP), each labeled (either radioactive label or fluorescent label with a different color of dye)

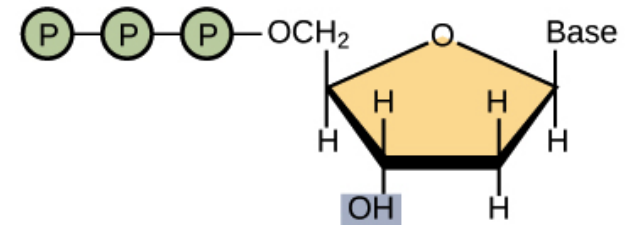


Dideoxy nucleotides are similar to regular, or deoxy, nucleotides, but with one key difference: **they lack a hydroxyl group on the 3' carbon of the sugar ring**. In a regular nucleotide, the 3' hydroxyl group acts as a "hook," allowing a new nucleotide to be added to an existing chain.

Once a dideoxy nucleotide has been added to the chain, there is no hydroxyl available and no further nucleotides can be added. The chain ends with the dideoxy nucleotide, which is marked with a particular color of dye depending on the base (A, T, C or G) that it carries.



Dideoxynucleotide (ddNTP)

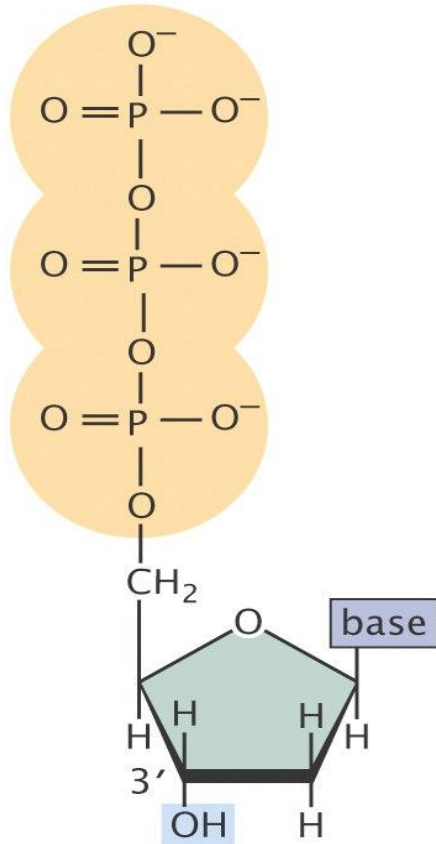


Deoxynucleotide (dNTP)

TP3: DNA Sequencing

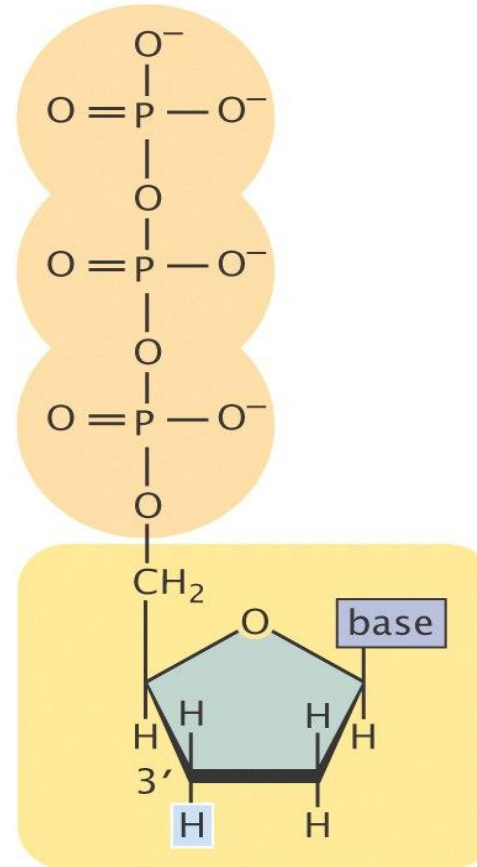
ddNTPs

A 3'-OH in normal DNA is necessary for elongation



Deoxyribonucleoside triphosphate (dNTP)

2'- deoxyribose



Dideoxyribonucleoside triphosphate (ddNTP)

2', 3'- dideoxyribose

The dideoxy sequencing requires a **special substrate** for DNA synthesis

dNTP vs ddNTP

Didesoxirribonucleosido trifosfato (ddNTP)

TP3: DNA Sequencing

ddNTPs labelling

- **Manual DNA sequencing**

Radioactive labeling

<https://www.youtube.com/watch?v=aPN8LP4YxPo>

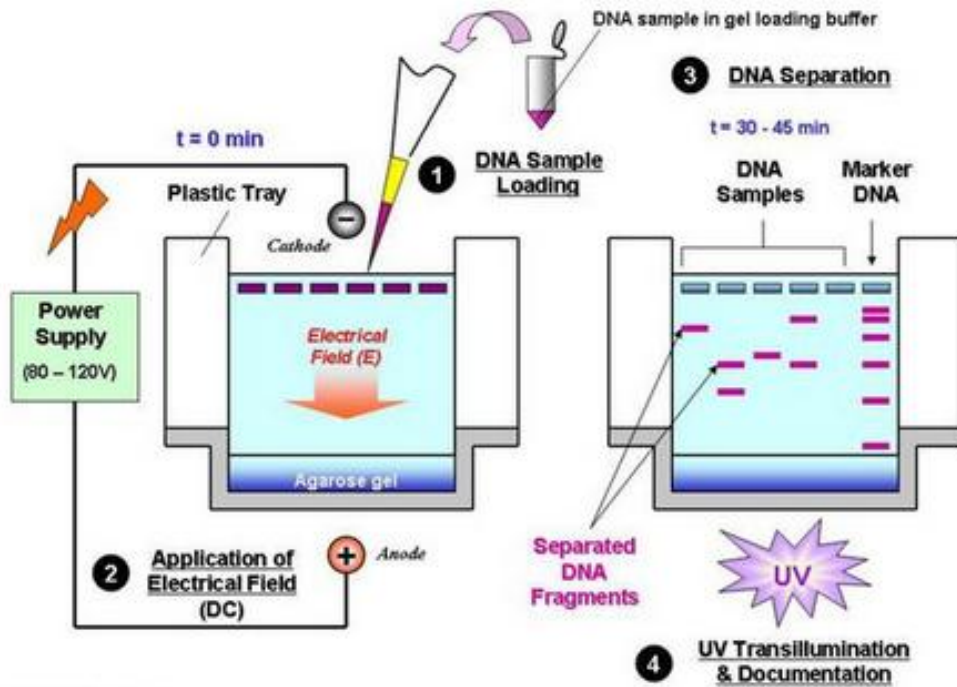
- **Automated DNA sequencing**

Fluorescence labeling with different fluorochromes

<https://www.youtube.com/watch?v=e2G5zx-OJlw>

TP3: DNA Sequencing

Manual DNA sequencing



Graphic©ESchmid/2001

Polyacrylamide gel electrophoresis separates ssDNA molecules that differ in length by just **one nucleotide**

Molecules are labelled with a radioactive protein or radioactive isotope, visualized by autoradiography producing a banding pattern

TP3: DNA Sequencing

Reading a sequencing gel

- You begin from the bottom where the smallest DNA fragments are,
- The sequence that you read will be in the 5'-3' direction,
- This sequence will be complementary to the template DNA chain

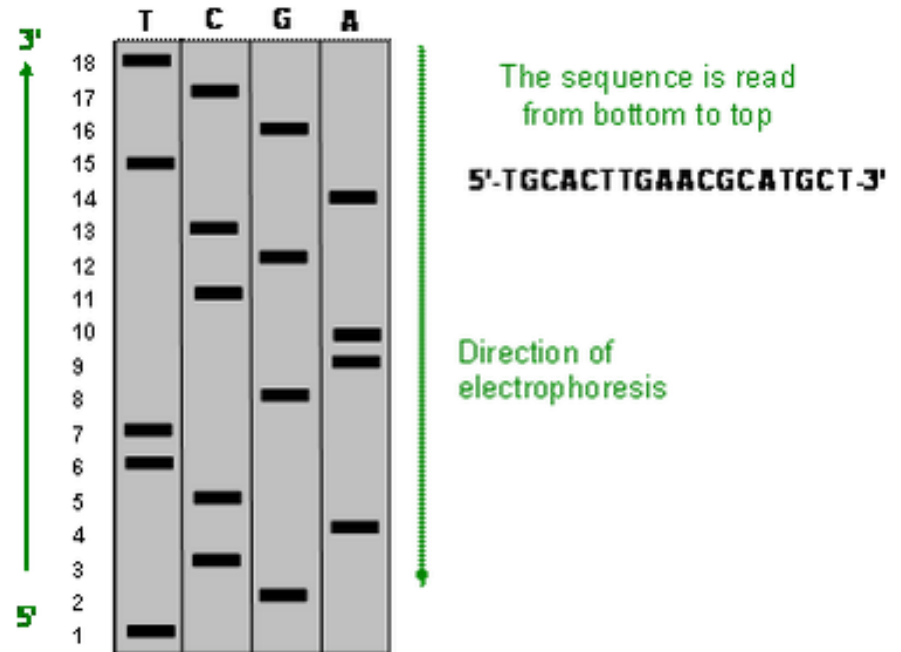
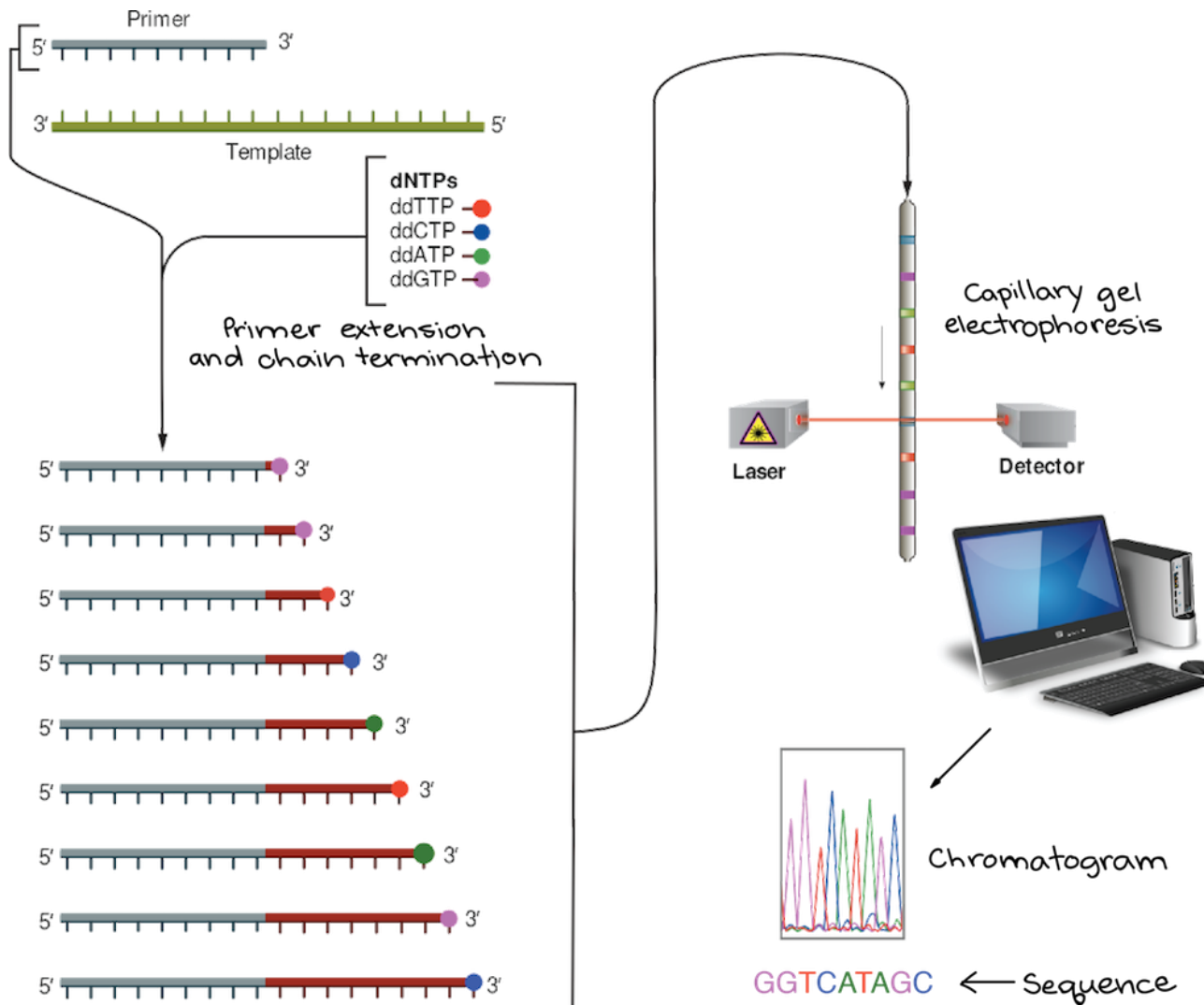


Figure 2. Reading the sequence from the gel.

TP3: DNA Sequencing

Automated DNA sequencing



TP3: DNA Sequencing

5' CCTATTATGACACAACCGCA 3'

ddCTP ddGTP ddTTP ddATP

C G T A

dNTPs

Template strand

Primer
(sequence known)

5' CCTATTATGACACAACCGCA 3'
3' GCGT 5'

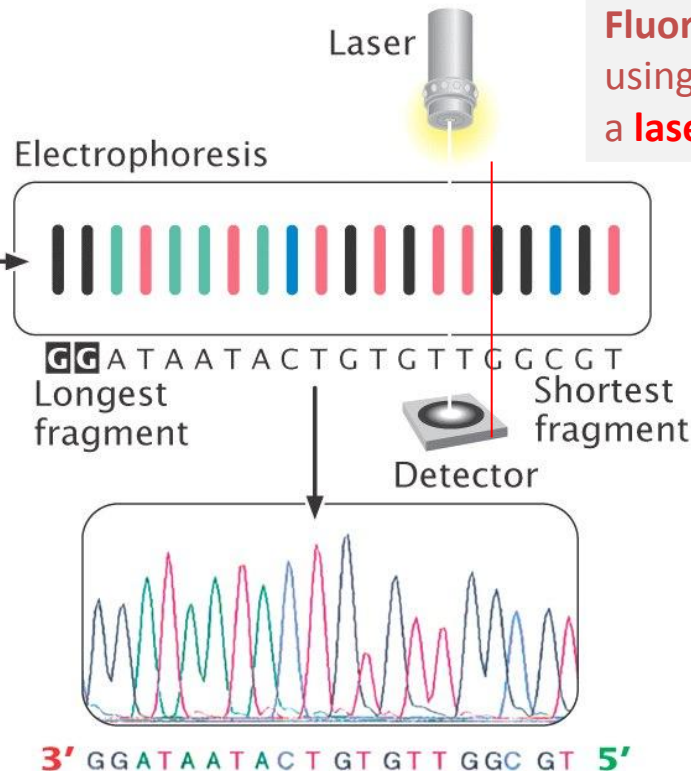
5' CCTATTATGACACAACCGCA 3'
3' GGATAA TACTGTGTTGCCGT 5'

5' CCTATTATGACACAACCGCA 3'
3' GATAA TACTGTGTTGCCGT 5'

Each of the four ddNTPs is tagged with a **fluorescent dye**

Fluorescent dye detected by using a **laser** beam and a **detector**

Denatured DNA products are mixed and loaded into a **single well** on an electrophoresis gel.



The sequence information is directly **read** and **electronically stored** into the computer, which converts it into the complementary-target-sequence

TP3: DNA Sequencing

Sequencing technology advances

- 1868: Discovery of DNA
- 1953: Watson and Crick propose double helix structure
- 1977: Sanger sequencing
- 1985: PCR
- 2000: Working draft human genome announced (Sanger method)
- 2005: 454 sequencer launch (pyrosequencing)
- 2006: Genome Analyzer launched (Solexa sequencing)
- 2007: SOLiD launched (ligation sequencing)
- 2009: Whole human genome no longer merits Nature/Science paper
- 2011: Illumina sequencer (sequencing by synthesis)
- 2011: Ion torrent
- 2011-18: 3rd generation sequencing: Pacbio, Oxford nanopore

\$ human
Genome

\$3 billion

\$2-3 million

\$250k

\$50k

\$20k

\$20k

\$20k

?<\$5k?



TP3: DNA Sequencing

Sequencing technology advances

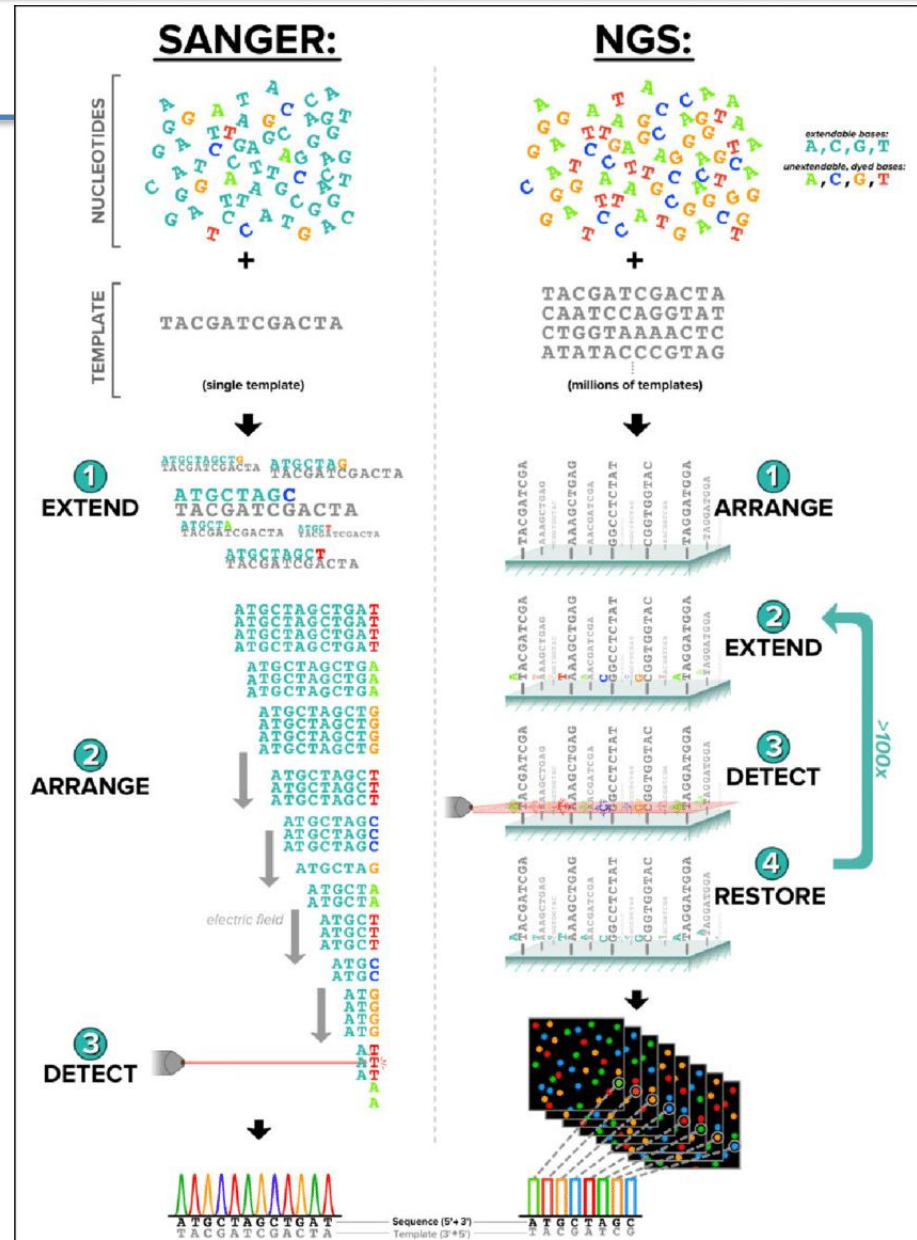
TABLE 1 | Comparison of available NGS platforms

Company	Read length	Applications	Website
454/Roche	400 bp (single end)	Bacterial and viral genomes, multiplex-PCR products, validation of point mutations, targeted somatic-mutation detection	http://www.454.com/
Illumina	150–300 bp (paired end)	Complex genomes (human, mouse and plants) and genome-wide NGS applications, RNA-seq, hybrid capture or multiplex-PCR products, somatic-mutation detection, forensics, noninvasive prenatal testing	http://www.illumina.com/
ABI SOLiD	75 bp (single end) or 50 bp (paired end)	Complex genomes (human, mouse, plants) and genome-wide NGS applications, RNA-seq, hybrid capture or multiplex-PCR products, somatic-mutation detection	http://www.thermofisher.com/us/en/home/life-science/sequencing/next-generation-sequencing/solid-next-generation-sequencing.html/
Pacific Biosciences	Up to 40 kb (single end or circular consensus)	Complex genomes (human, mouse and plants), microbiology and infectious-disease genomes, transcript-fusion detection, methylation detection	http://www.pacb.com/
Ion Torrent	200–400 bp (single end)	Multiplex-PCR products, microbiology and infectious diseases, somatic-mutation detection, validation of point mutations	http://www.thermofisher.com/us/en/home/life-science/sequencing/next-generation-sequencing.html/
Oxford Nanopore	Variable: depends on library preparation (1D or 2D reads)	Pathogen surveillance, targeted mutation detection, metagenomics, bacterial and viral genomes	http://nanoporetech.com/
Qiagen GeneReader	107 bp (single end)	Targeted mutation detection, liquid biopsy in cancer	http://www.genereaderngs.com/

TP3: DNA Sequencing

Sanger sequencing vs NGS

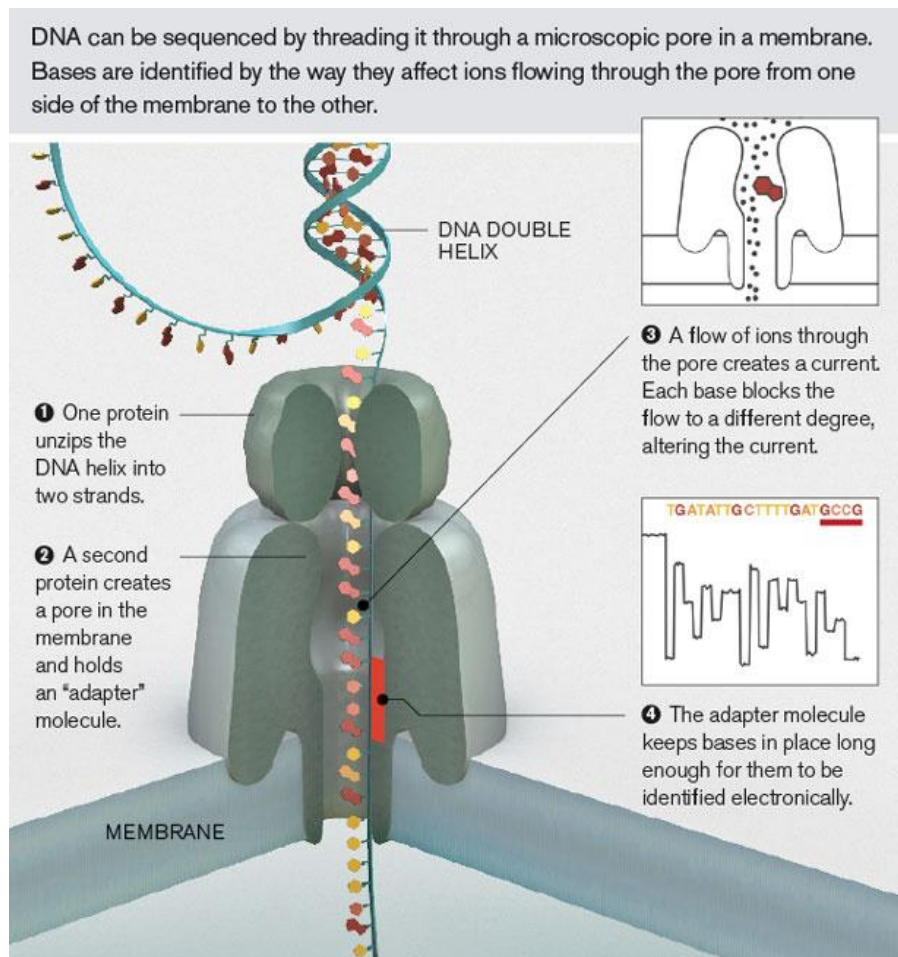
<https://www.thermofisher.com/blog/behindthebench/when-do-i-use-sanger-sequencing-vs-ngs-seq-it-out-7/>



TP3: DNA Sequencing

NGS latest developments

Nanopore sequencing:

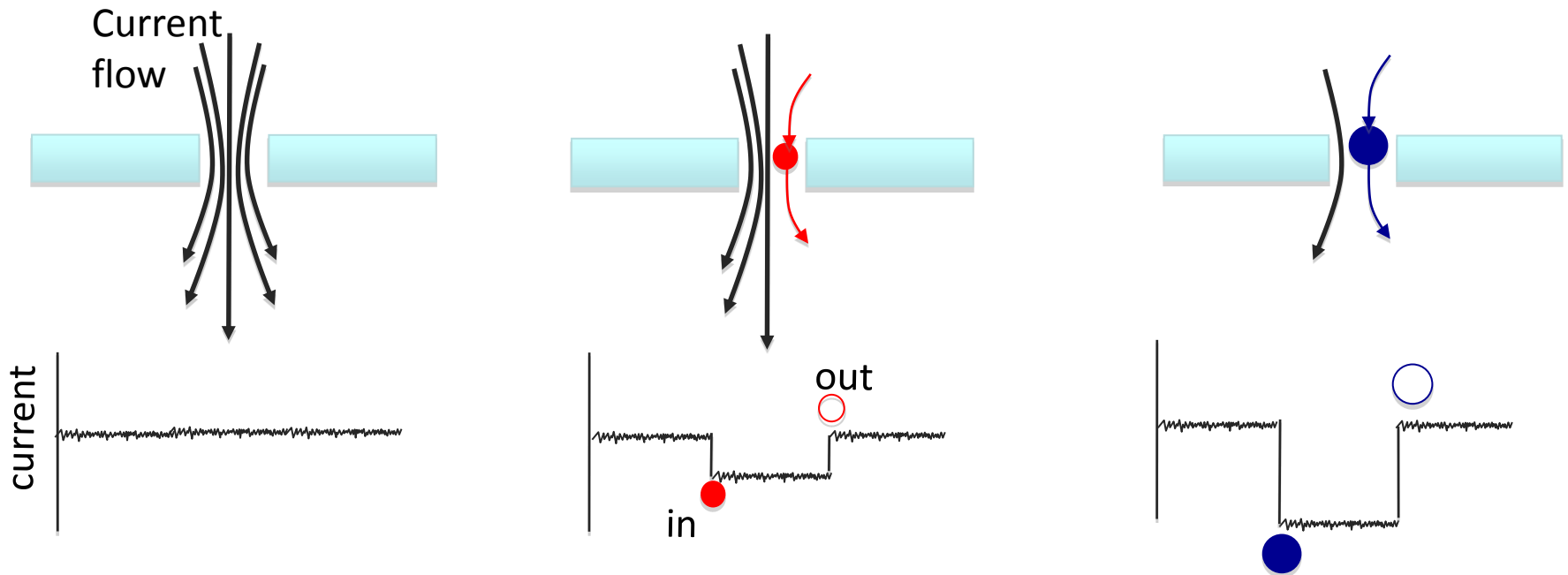


Determine the sequence of DNA fragments by passing DNA through a protein (or other) pore in a membrane

TP3: DNA Sequencing

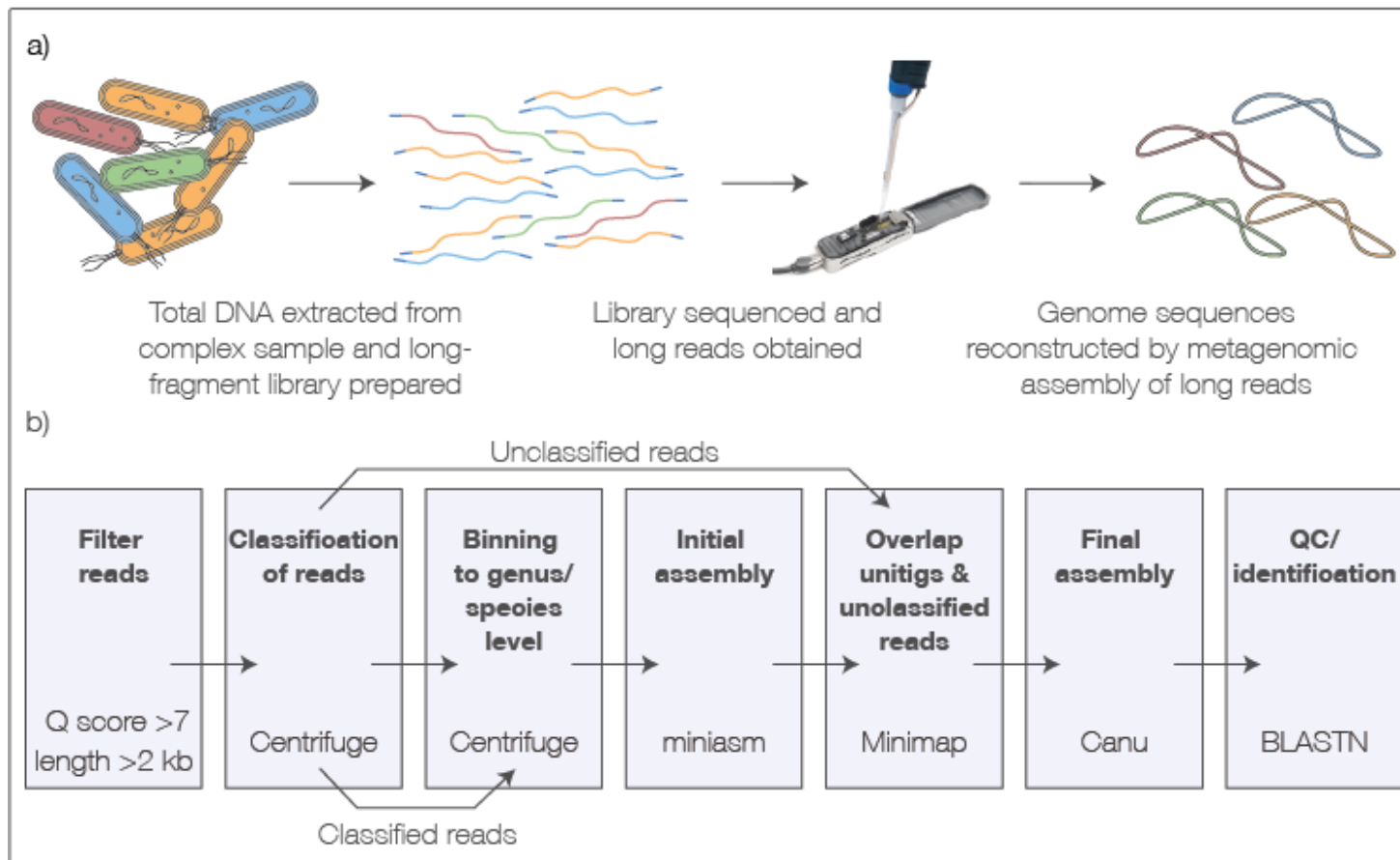
NGS latest developments

- Nanopore = ‘very small hole’
- Electrical current flows through the hole
- Introduce analyte of interest into the hole → identify “analyte” by the disruption or block to the electrical current



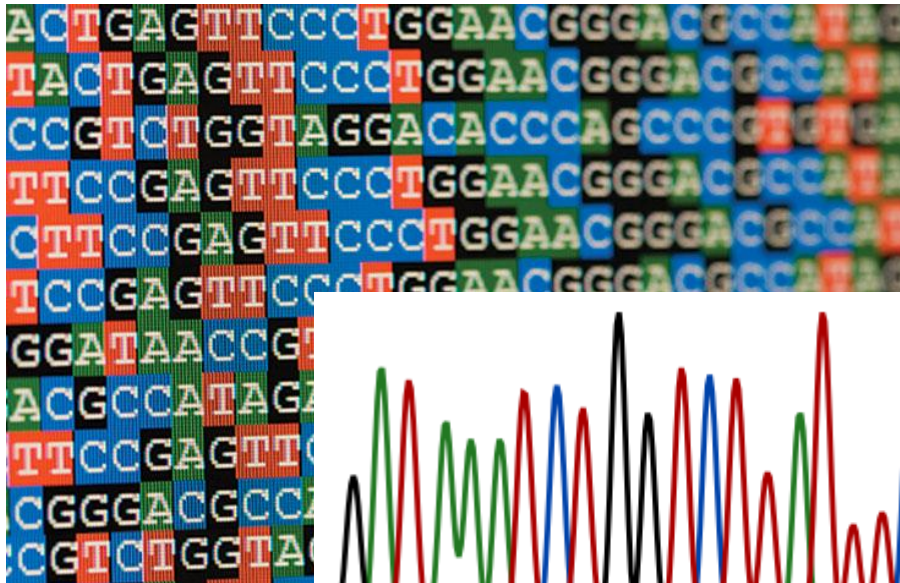
TP3: DNA Sequencing

NGS latest developments



TP3: DNA Sequencing

Sequencing data analysis



120 130
G A T A A A T C T G G T C T T A T T C C



Databases

Sequence alignment of different clones or reads

Nucleotide or peptide sequence comparison with other species (blast)

Sequence analysis for:

Genome comparisons

Restriction map

ORFs

Peptidic sequence

Specific sequences (promoter, DNA-binding domains (ex. response elements), *stem-loop*, palindrom, direct and inverted repeats etc)

% G/C

Codon usage (codon preference)

TP3: DNA Sequencing

Sequencing data analysis

Major Sequence Repositories

GenBank or NCBI (all known nucleotide and protein sequences)

www.ncbi.nlm.nih.gov/Web/Genbank/

Ensembl (all known nucleotide and protein sequences)

www.ensembl.org/index.html

Genome Databases

Flybase (Drosophyla sequences and genomic information)

www.fruitfly.org

MGD (Mouse genetics and genomics)

www.informatics.jax.org

Grapevine

<http://genomes.cribi.unipd.it/grape/>

Arabidopsis

<https://www.arabidopsis.org/>

Genetic Maps

GBD (Human genes and genomic maps)

www.gbd.org

NCBI genome mapping

<https://www.ncbi.nlm.nih.gov/probe/docs/applmapping/>

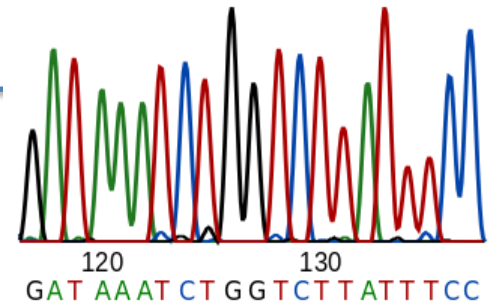
Gene Expression

BodyMap (Human and mouse gene expression data)

bodymap.ims.u-tokyo.ac.jp

Tair

OPANDA



Gene Identification and Structure

EID (Protein-coding, intron-containing genes)

mcb.harvard.edu/gilbert/EID/

Exint (Exon-intron structure of eukaryotic genes)

intron.bic.nus.edu.sg/exint/extint.html

TRRD (Regulatory regions of eukaryotic genes)

www.mgs.bionet.nsc.re/mgs/dbases/trrd4/

Protein interaction database

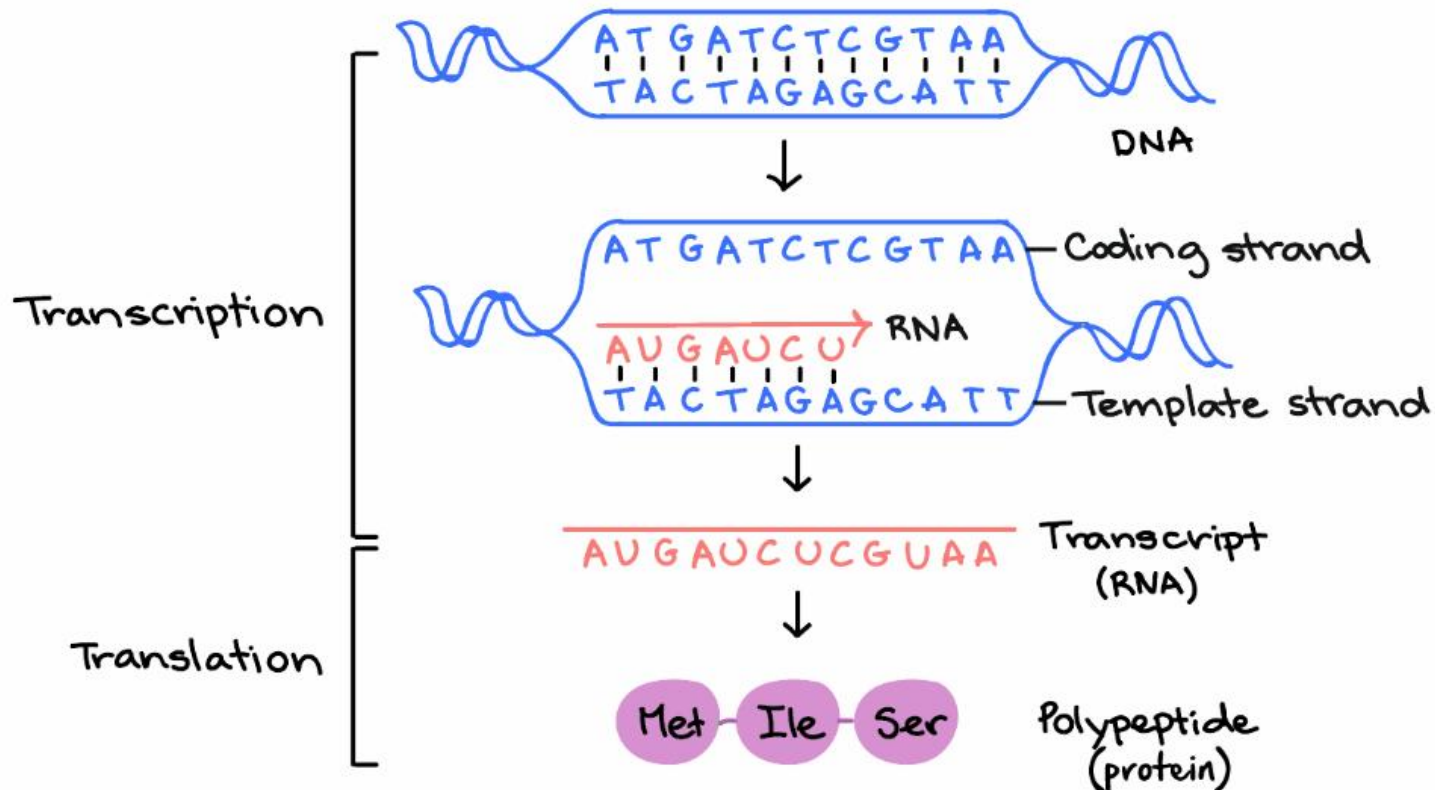
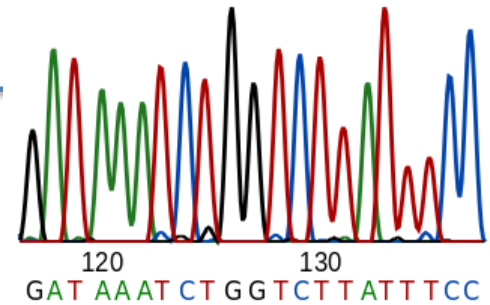
String

<https://string-db.org/>

TP3: DNA Sequencing

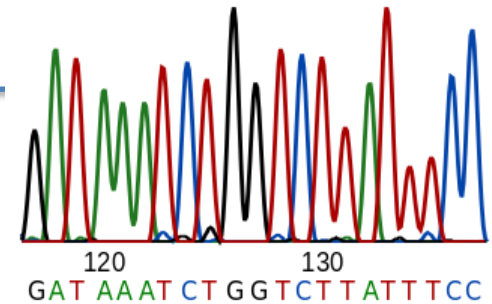
Sequencing data analysis - ORF

Definition of the **open reading frame**: (ORF) is the part of a reading frame that has the potential to code for a protein or peptide. An ORF is a continuous stretch of codons beginning with a start codon (usually **AUG**) and ending with a stop codon (usually **TAA**, **TAG** or **TGA**)



TP3: DNA Sequencing

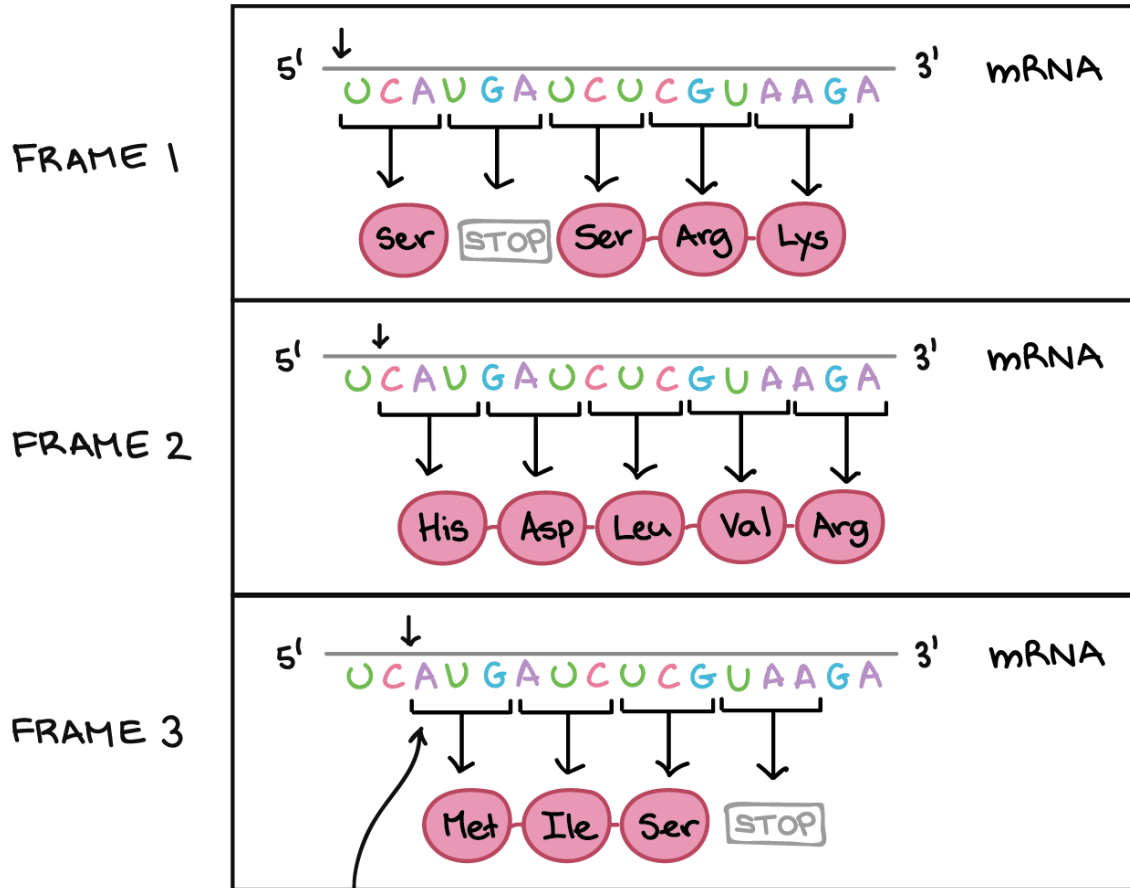
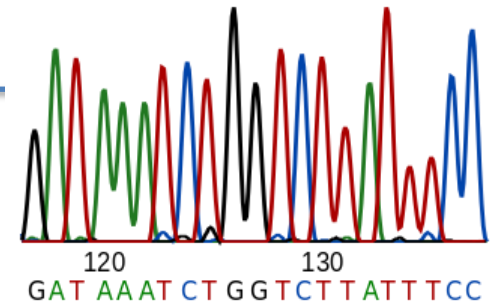
Sequencing data analysis - ORF



5'-Base		Middle	Base		3'-Base
	U(=T)	C	A	G	
U(=T)	Phe	Ser	Tyr	Cys	U(=T)
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	Term	Term	A
	Leu	Ser	Term	Trp	G
C	Leu	Pro	His	Arg	U(=T)
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U(=T)
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U(=T)
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

TP3: DNA Sequencing

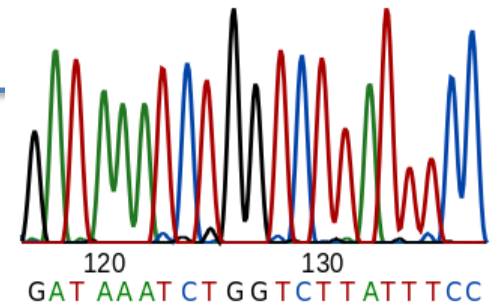
Sequencing data analysis - ORF



Start codon's position ensures that this frame is chosen

TP3: DNA Sequencing

Sequencing data analysis - ORF



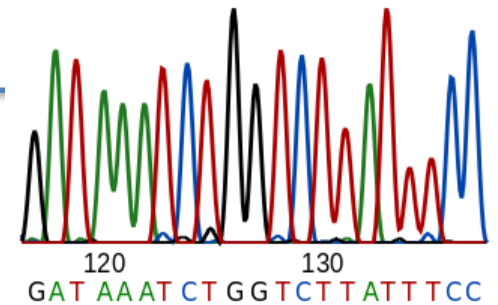
5' AGGTGACACCGCAAGCCTTATATTAGC 3'

Illustration of possible reading frames:

AGG·TGA·CAC·CGC·AAG·CCT·TAT·ATT·AGC
A·GGT·GAC·ACC·GCA·AGC·CTT·ATA·TTA·GC
AG·GTG·ACA·CCG·CAA·GCC·TTA·TAT·TAG·C

TP3: DNA Sequencing

Sequencing data analysis - ORF



5' atgccc aagctgaatagcgtagagggggttttcatcatttgaggacgatgtataa 3'

1	atg	ccc	aag	ctg	aat	agc	gta	gag	ggg	ttt	tca	tca	ttt	gag	gac	gat	gta	taa
	M	P	K	L	N	S	V	E	G	F	S	S	F	E	D	D	V	*
2	tgc	cca	agc	tga	ata	gcg	tag	agg	ggt	ttt	cat	cat	ttg	agg	acg	atg	tat	
	C	P	S	*	I	A	*	R	G	F	H	H	L	R	T	M	Y	
3	gcc	caa	gct	gaa	tag	cgt	aga	ggg	ggt	ttc	atc	att	tga	gga	cga	tgt	ata	
	A	Q	A	E	*	R	R	G	V	F	I	I	*	G	R	C	I	

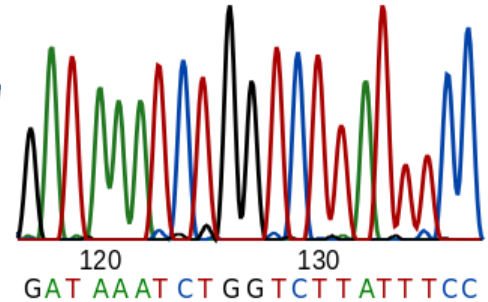
TP3: DNA Sequencing

Restriction map



NEBcutter V2.0

[Program Guide](#) [Help](#) [Comments](#)



This tool will take a DNA sequence and find the large, non-overlapping open reading frames using the E.coli genetic code and the sites for all Type II and commercially available Type III restriction enzymes that cut the sequence just once. By default, only enzymes available from NEB are used, but other sets may be chosen. Just enter your sequence and "submit". Further options will appear with the output. **The maximum size of the input file is 1 MByte, and the maximum sequence length is 300 KBases.**
[What's new in V2.0](#) [Citing NEBcutter](#)

Local sequence file: Nenhum ficheiro selecionado. Standard sequences:
GenBank number: [Browse GenBank](#) # Plasmid vectors
or paste in your DNA sequence: (plain or FASTA format) # Viral + phage

The sequence is: Linear Circular Enzymes to use: NEB enzymes
 All commercial enzymes
 All specific enzymes
 All + defined enzymes
 Only defined enzymes
[\[define oligo\]](#)

Minimum ORF length to display: a.a.

BamHI
1 gga tcc GCA GCG GAA ATC AGT GGT CAC ATC GTA CGT TCC CCG ATG GTT GGT ACT
1 Gly Ser Ala Ala Glu Ile Ser Gly His Ile Val Arg Ser Pro Met Val Gly Thr

SplI
55 TTC TAC CGC ACC CCA AGC CCG GAC GCA AAA GCT TTC ATC GAA GTG GGT CAG AAA
19 Phe Tyr Arg Thr Pro Ser Pro Asp Ala Lys Ala Phe Ile Glu Val Gly Gln Lys

HindIII
109 GTC AAC GTG GGC GAT ACC CTA TGC ATC GTT GAA GCC ATG AAA ATG ATG AAC CAG
37 Val Asn Val Gly Asp Thr Leu Cys Ile Val Glu Ala Met Met Lys Met Asn Gln

NsiI
163 ATC GAA GCG GAC AAA TCC GGT ACC GTG AAA GCA ATT CTG GTC GAA TCC GGA CAA
55 Ile Glu Ala Asp Lys Ser Gly Thr Val Lys Ala Ile Leu Val Glu Ser Gly Gln

KpnI **BspEI**
217 CCG GTA GAA TTT GAC GAG CCG CTG GTC GTC ATC GAG TAA gaa ttc
73 Pro Val Glu Phe Asp Glu Pro Leu Val Val Ile Glu ***

EcoRI